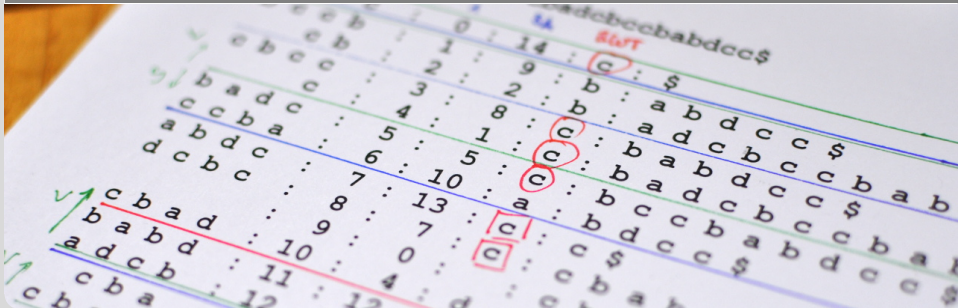


# Scalable Construction of Text Indexes with Thrill

Timo Bingmann, Simon Gog, and Florian Kurpicz · IEEE Big Data · December 12th, 2018

INSTITUTE OF THEORETICAL INFORMATICS – ALGORITHMICS



0 1 2 3 4 5 6 7 8 9 10 11 12 13

**Example  $T = [t\ o\ b\ e\ o\ r\ n\ o\ t\ t\ o\ b\ e\ \$]$**

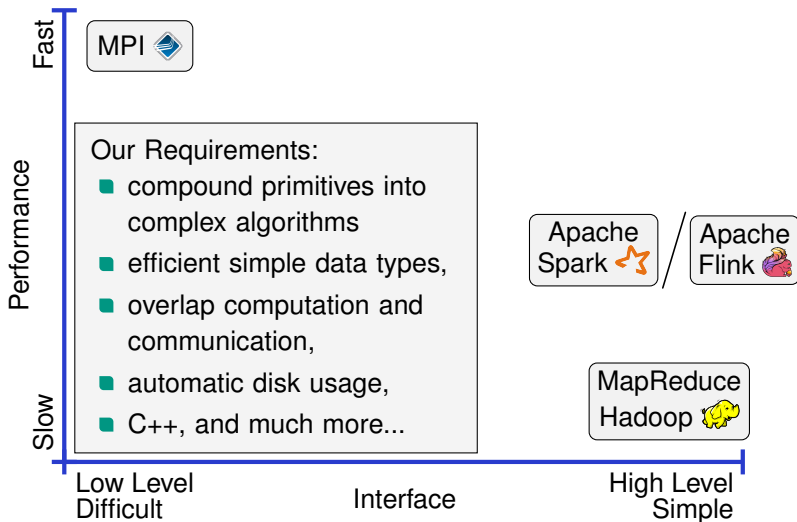
$i$	$T_i$
0	t o b e o r n o t t o b e \$
1	o b e o r n o t t o b e \$
2	b e o r n o t t o b e \$
3	e o r n o t t o b e \$
4	o r n o t t o b e \$
5	r n o t t o b e \$
6	n o t t o b e \$
7	o t t o b e \$
8	t t o b e \$
9	t o b e \$
10	o b e \$
11	b e \$
12	e \$
13	\$

Example  $T = [t\ o\ b\ e\ o\ r\ n\ o\ t\ t\ o\ b\ e\ \$]$

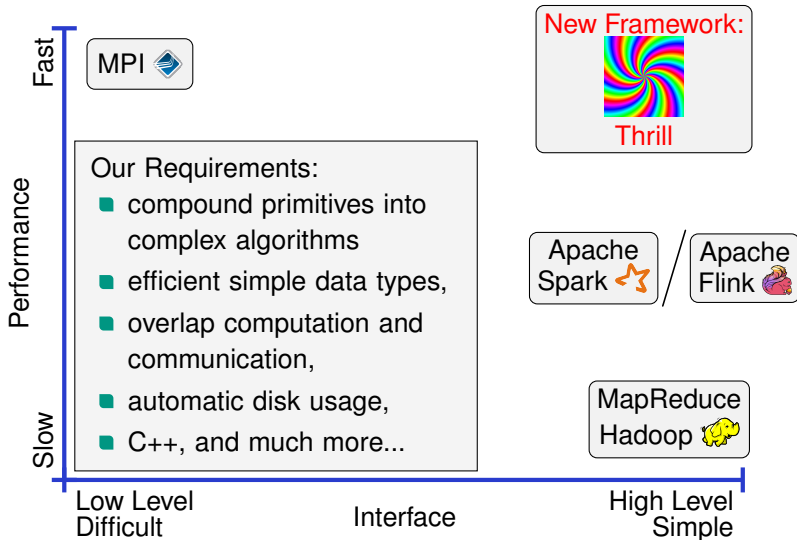
$SA_i$	$LCP_i$	$T_{SA_i \dots n}$
13	-	\$
11	0	b e \$
2	2	b e o r n o t t o b e \$
12	0	e \$
3	1	e o r n o t t o b e \$
6	0	n o t t o b e \$
10	0	o b e \$
1	3	o b e o r n o t t o b e \$
4	1	o r n o t t o b e \$
7	1	o t t o b e \$
5	0	r n o t t o b e \$
9	1	t o b e \$
0	4	t o b e o r n o t t o b e \$
8	1	t t o b e \$



# Big Data Batch Processing



# Big Data Batch Processing



# Thrill's Goal and Current Status

An **easy way** to program **fast distributed** algorithms in C++.

## Current Status:

- Open-source prototype at <http://github.com/thrill/thrill>.
- $\approx$  60 K lines of C++14 code, written by  $\geq$  12 contributors.
- Published at [IEEE Conference on Big Data](#) [B, et al. '16]
- Faster than Apache Spark and Apache Flink on [five micro benchmarks](#): WordCount1000, WordCountCC, PageRank, TeraSort, and K-Means.

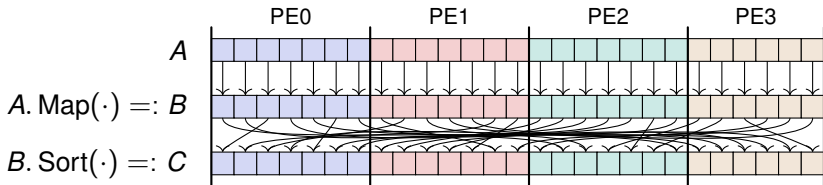
## Case Studies:

- **Five suffix sorting algorithms** [B, Gog, Kurpicz, [this presentation](#)]
- Louvain graph clustering algorithm [Hamann et al. [Euro-Par'18](#)]
- Process scientific data on HPC (poster) [Karabin et al. [SC'18](#)]
- More examples: stochastic gradient descent, triangle counting, etc.
- **Future**: fault tolerance, scalability, and more applications.

# Distributed Immutable Array (DIA)

## User Programmer's View:

- $\text{DIA}\langle T \rangle$  = distributed array of items  $T$  on the cluster
- Cannot access items directly, instead use small set of **scalable primitives**, for example: **Map**, **Sort**, **ReduceByKey**, **Zip**, **Window**, etc.

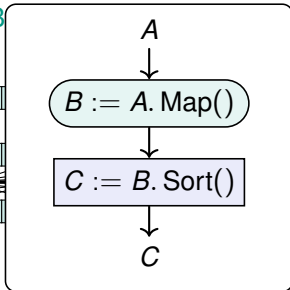
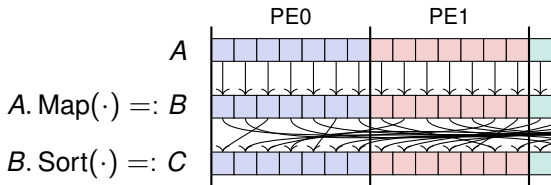




# Distributed Immutable Array (DIA)

## User Programmer's View:

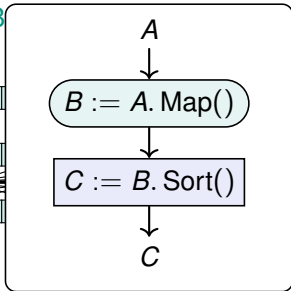
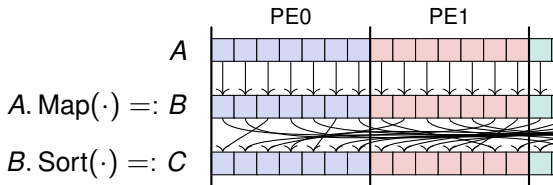
- $\text{DIA}\langle T \rangle =$  distributed array of items  $T$  on the cluster
- Cannot access items directly, instead use small set of scalable primitives, for example: **Map**, **Sort**, **ReduceB**



# Distributed Immutable Array (DIA)

## User Programmer's View:

- $\text{DIA}\langle T \rangle$  = distributed array of items  $T$  on the cluster
- Cannot access items directly, instead use small set of scalable primitives, for example: **Map**, **Sort**, **ReduceB**



## Framework Designer's View:

- Goals: distribute work, optimize execution on cluster, add redundancy where applicable.  $\Rightarrow$  build data-flow graph.
- $\text{DIA}\langle T \rangle$  = pipelined chain of computations

# Distributed Suffix Sorting

## Prior Work on **Distributed** Algorithms:

- Early research with distributed sorting algorithms and using “pruned suffixes” and “**sistrings**” with additional characters. [KRRNZ, ICAPP’97]
- Emulation of sequential suffix sorting algorithms. [KN, SPIRE’99]
- With assumption that the **entire** text is on each processor. [FAK, ’01]
- First MPI implementation of **DC3** using sample sort. [KS, PVM/MPI’06]
- Hadoop implementation of simple algorithms. [MBS, MapReduce’11]
- Implementation of [FAK01] named cloudSACA for Amazon Web Services. [AKA, MECBME’14]
- Implementation of **prefix doubling** using MPI. [FA, SC’15]

# Paper: Scalable Suffix Sorting with Thrill

Implemented **five algorithms** for **distributed external** memory

- Prefix doubling using the inverse suffix array (DoublingW)
- Prefix doubling using sorting (DoublingS)
- Prefix doubling with discarding (Discarding)
- DC3
- DC7



Example  $T = [t\ o\ b\ e\ o\ r\ n\ o\ t\ t\ o\ b\ e\ \$]$

$i$	$T_i$
0	t o b e o r n o t t o b e \$
1	o b e o r n o t t o b e \$
2	b e o r n o t t o b e \$
3	e o r n o t t o b e \$
4	o r n o t t o b e \$
5	r n o t t o b e \$
6	n o t t o b e \$
7	o t t o b e \$
8	t t o b e \$
9	t o b e \$
10	o b e \$
11	b e \$
12	e \$
13	\$

Example  $T = [t\ o\ b\ e\ o\ r\ n\ o\ t\ t\ o\ b\ e\ \$]$

$i$	$T_i$
13	0 \$
2	1 b e o r n o t t o b e \$
11	
3	3 e o r n o t t o b e \$
12	
6	5 n o t t o b e \$
1	
4	6 o b e o r n o t t o b e \$
7	
10	
5	10 r n o t t o b e \$
0	
8	11 t o b e o r n o t t o b e \$
9	

Example  $T = [t\ o\ b\ e\ o\ r\ n\ o\ t\ t\ o\ b\ e\ \$]$

$i$	$T_i$
13	0 \$
2	1 b e <sup>3</sup> o r n o t t o b e \$
11	1 b e <sup>3</sup> \$
3	3 e o <sup>6</sup> r n o t t o b e \$
12	3 e \$ <sup>0</sup>
6	5 n o t t o b e \$
1	6 o b <sup>1</sup> e o r n o t t o b e \$
4	6 o r <sup>10</sup> o t t o b e \$
7	6 o t <sup>11</sup> o b e \$
10	6 o b <sup>1</sup> e \$
5	10 r n o t t o b e \$
0	11 t o <sup>6</sup> b e o r n o t t o b e \$
8	11 t t <sup>11</sup> o b e \$
9	11 t o <sup>6</sup> b e \$

Example  $T = [t\ o\ b\ e\ o\ r\ n\ o\ t\ t\ o\ b\ e\ \$]$

$i$	$T_i$
13	0 \$
2	1 b e 3 o r n o t t o b e \$
11	1 b e 3 \$
12	3 e \$ 0
3	5 e o 6 r n o t t o b e \$
6	5 n o t t o b e \$
1	6 o b 1 e o r n o t t o b e \$
10	6 o b 1 e \$
4	10 o r 10 o t t o b e \$
7	11 o t 11 o b e \$
5	10 r n o t t o b e \$
0	11 t o 6 b e o r n o t t o b e \$
9	11 t o 6 b e \$
8	11 t t 11 o b e \$



Example  $T = [t\ o\ b\ e\ o\ r\ n\ o\ t\ t\ o\ b\ e\ \$]$

$i$	$T_i$
13	0 \$
2	1 b e o8r n o t t o b e \$
11	b e \$0
12	3 e \$
3	4 e o r n o t t o b e \$
6	5 n o t t o b e \$
1	6 o b e4o r n o t t o b e \$
10	o b e3\$
4	8 o r n o t t o b e \$
7	9 o t t o b e \$
5	10 r n o t t o b e \$
0	11 t o b1e o r n o t t o b e \$
9	t o b1e \$
8	13 t t o b e \$

Example  $T = [t\ o\ b\ e\ o\ r\ n\ o\ t\ t\ o\ b\ e\ \$]$

$i$	$T_i$
13	0 \$
11	1 b e \$0
2	2 b e o8r n o t t o b e \$
12	3 e \$
3	4 e o r n o t t o b e \$
6	5 n o t t o b e \$
10	6 o b e \$
1	7 o b e o r n o t t o b e \$
4	8 o r n o t t o b e \$
7	9 o t t o b e \$
5	10 r n o t t o b e \$
0	11 t o b1e o r n o t t o b e \$
9	t o b1e \$
8	13 t t o b e \$

Example  $T = [t\ o\ b\ e\ o\ r\ n\ o\ t\ t\ o\ b\ e\ \$]$

$i$	$T_i$
13	0 \$
11	1 b e \$
2	2 b e o r n o t t o b e \$
12	3 e \$
3	4 e o r n o t t o b e \$
6	5 n o t t o b e \$
10	6 o b e \$
1	7 o b e o r n o t t o b e \$
4	8 o r n o t t o b e \$
7	9 o t t o b e \$
5	10 r n o t t o b e \$
0	t o b e o r n o t t o b e \$
9	11 t o b e \$ 0
8	13 t t o b e \$

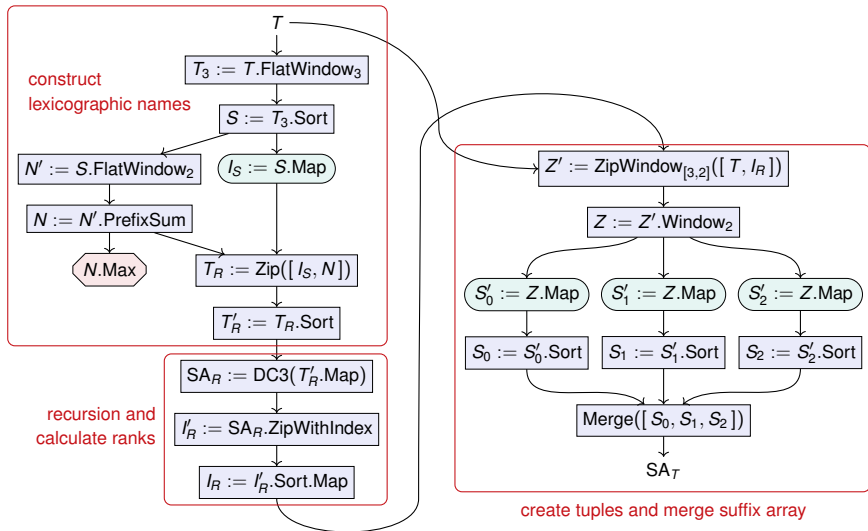
Example  $T = [t\ o\ b\ e\ o\ r\ n\ o\ t\ t\ o\ b\ e\ \$]$

$i$	$T_i$
13	0 \$
11	1 b e \$
2	2 b e o r n o t t o b e \$
12	3 e \$
3	4 e o r n o t t o b e \$
6	5 n o t t o b e \$
10	6 o b e \$
1	7 o b e o r n o t t o b e \$
4	8 o r n o t t o b e \$
7	9 o t t o b e \$
5	10 r n o t t o b e \$
9	11 t o b e \$ 0
0	12 t o b e o r 8 n o t t o b e \$
8	13 t t o b e \$

# Thrill Pseudo-Code for Prefix Doublings

```
1 Function PrefixDoublingWithSorting( $T \in \langle \Sigma \rangle$ )
2    $S := T.$ Window2(( $i, [t_0, t_1]$ )  $\mapsto$  ( $i, t_0, t_1$ ))    // Initial triples ( $i, T[i], T[i + 1]$ ).
3   for  $k := 1$  to  $\lceil \log_2 |T| \rceil - 1$  do
4      $S := S.$ Sort(( $i, n_0, n_1$ ) by ( $n_0, n_1$ ))           // Sort triples by name pair.
5      $N := S.$ FlatWindow2(( $j, [a, b]$ )  $\mapsto$  CmpName( $j, a, b$ )) // Outputs 0 or  $j$ .
6     if  $N.$ Filter(( $i, n$ )  $\mapsto$  ( $n = 0$ )).Size() = 1 then // If all names distinct, then
7       return  $N.$ Map(( $i, n$ )  $\mapsto$   $i$ )                    // return names as suffix array, else
8      $N := N.$ PrefixSum(( $i, n$ ), ( $i', n'$ )  $\mapsto$  ( $i', \max(n, n')$ )) // make new names.
9      $N := N.$ Sort(( $i, n$ ) by  $i$ )                          // Compute ISA2k.
10     $S := N.$ Window2k+1(( $j, [(i, n), \dots, (i', n')]$ )  $\mapsto$ 
       $\left\{ \begin{array}{ll} (i, n, n') & \text{if } j + 2^k < |T|, \\ (i, n, 0) & \text{otherwise.} \end{array} \right\}$  // Compare names ISA2k[ $i$ ]
      // and ISA2k[ $j + 2^k$ ].
```

# Data-Flow Graph of DC3 with Recursion



# Experiments on AWS EC2

**Inputs:** prefixes of

- **Wikipedia** XML dump (up to 125.6 GiB)
- **Gutenberg** text document corpus (up to 23 GiB)
- **Pi** digits of  $\pi$  (“3.1415926535...”)

**Machine:**

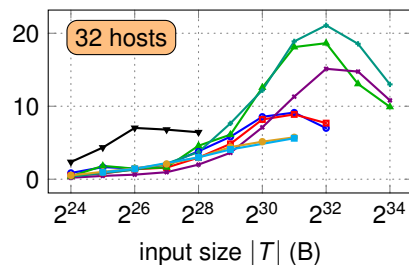
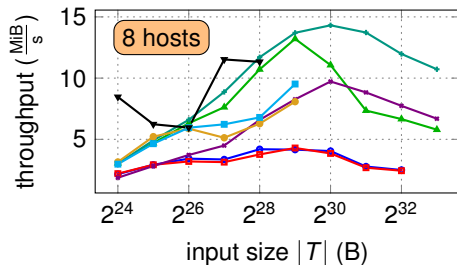
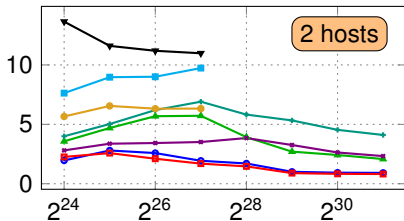
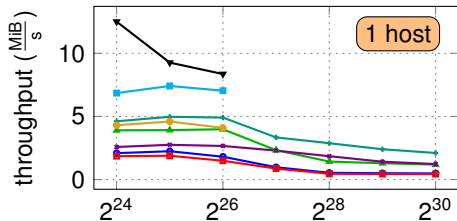
- up to  $32 \times$  **i3.4xlarge** EC2 instances.
- **16-core** Intel Xeon E5-2686 CPUs with 2.30 GHz
- **8 GB** of RAM, and  $2 \times$  **1.9 TB** NVMe SSDs



**Competitors:**

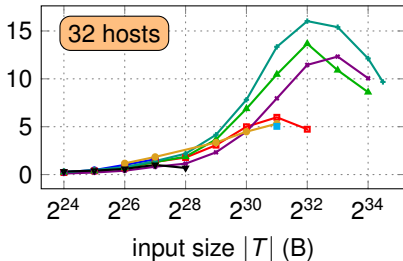
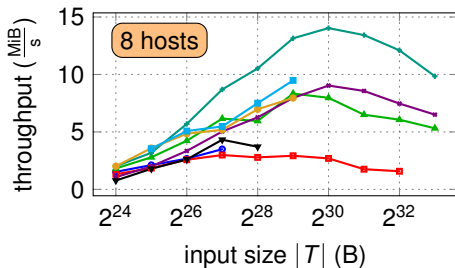
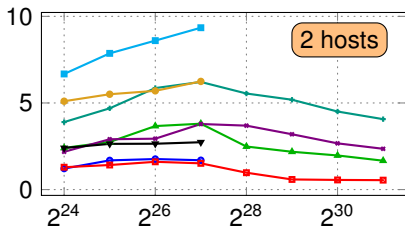
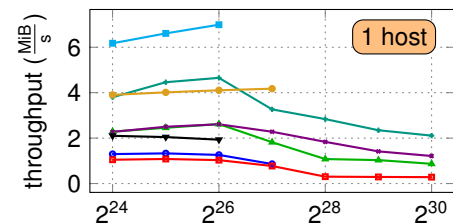
- **BKS.pDC3** and **BKS.pDC7** – MPI difference cover algorithm
- **FA.psac** – MPI prefix doubling algorithm
- also compared against fastest non-distributed algorithms:  
**M.divsufsort**, **M.divsufsort.par**, and **M.sais**.

# Suffix Sorting Wikipedia on AWS EC2



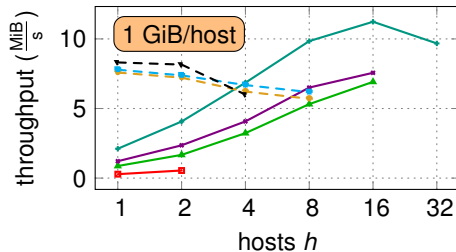
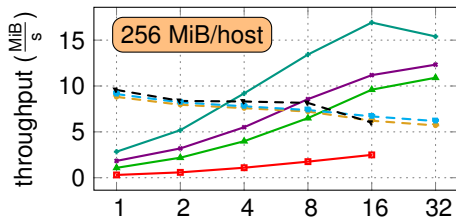


# Suffix Sorting Gutenberg on AWS EC2

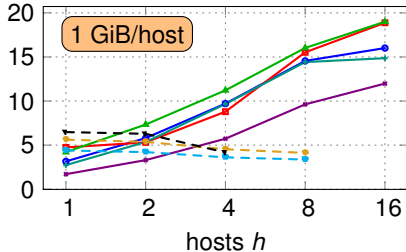
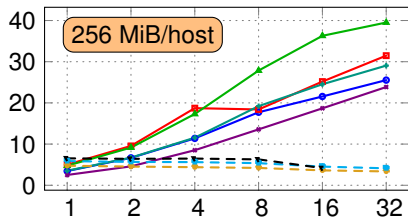


# Weak Scaling and COST on AWS EC2

## Gutenberg

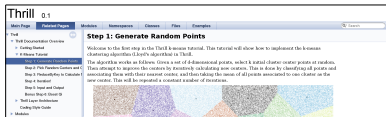


## Pi



# Current and Future Work

- Open-Source at <http://project-thrill.org> and Github.
- High quality, **very modern C++14** code.
- A **K-Mean tutorial** is available!



## Ideas for Future Work:

- Distributed rank()/select() and succinct bit vectors for text search?
- Beyond DIA<T>? Graph<V,E>? Matrix<T>?
- **Fault tolerance** in the algorithms and **scalability** to large clusters.

Thank you for your attention!  
Questions?